# STAR *fileCatalog, tagDB*
## and Grand Challenge Architecture

A. Vaniachine

**presenting for the Grand Challenge Collaboration**

**(http:/www-rnc.lbl.gov/GC/)**

**April 9, 2000**

**Joint ALICE-STAR Computing Meeting**

# Outline

- **GCA Overview**

- **STAR Interface:**
  - *fileCatalog*
  - *tagDB*
  - *StChallenger*

- **Current Status**

- **Conclusion**

# GCA: Grand Challenge Architecture

- *An order-optimized prefetch architecture for data retrieval from multilevel storage in a multiuser environment*

- **Queries select events and specific event components based upon tag attribute ranges**
  - query estimates are provided prior to execution
  - collections as queries are also supported

- **Because event components are distributed over several files, processing an event requires delivery of a "bundle" of files**

- **Events are delivered in an order that takes advantage of what is already on disk, and multiuser policy-based prefetching of further data from tertiary storage**

- **GCA intercomponent communication is CORBA-based, but physicists are shielded from this layer**

# Participants

- **NERSC/Berkeley Lab**
  - L. Bernardo, A. Mueller, H. Nordberg, A.Shoshani, A. Sim, J. Wu
- **Argonne**
  - D. Malon, E. May, G. Pandola
- **Brookhaven Lab**
  - B. Gibbard, S. Johnson, J. Porter, T.Wenaus
- **Nuclear Science/Berkeley Lab**
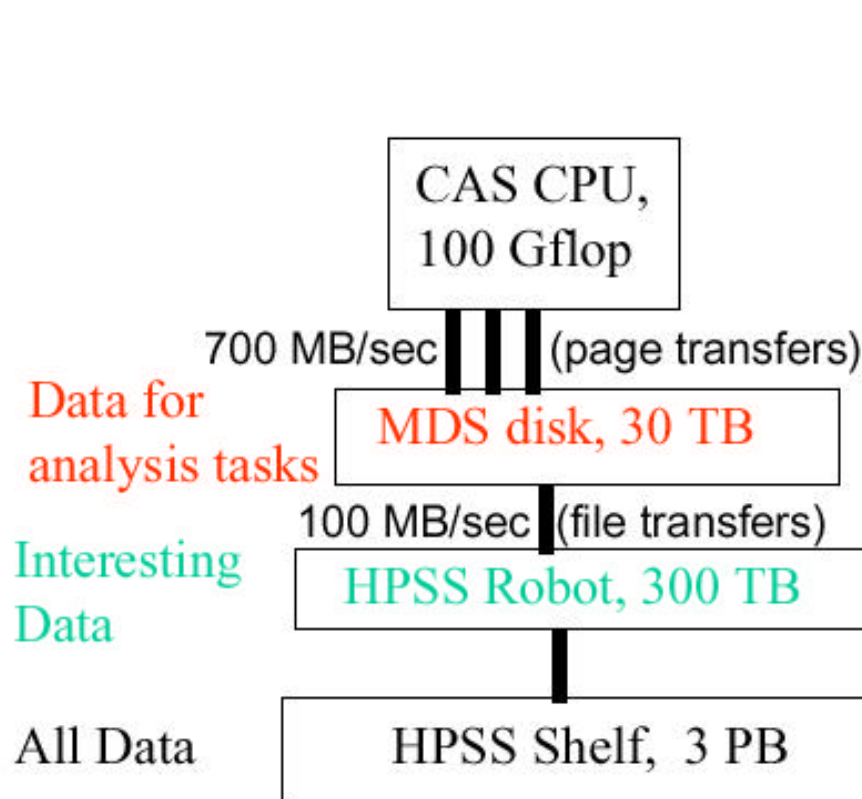  - D. Olson, A. Vaniachine, J. Yang, D.Zimmerman

# Problem

- **There are several**
  - **Not all data fits on disk ($$)**
    - **Part of 1 year's DST's fit on disk**
      - **What about last year, 2 year's ago?**
      - **What about hits, raw?**
  - **Available disk bandwidth means data read into memory must be efficiently used ($$)**
    - **don't read unused portions of the event**
    - **Don't read events you don't need**
  - **Available tape bandwidth means files read from tape must be shared by many users, files should not contain unused bytes ($$$$)**
  - **Facility resources are sufficient only if used efficiently**
    - **Should operate steady-state (nearly) fully loaded**

# Bottleneks



CAS CPU, 100 Gflop

700 MB/sec (page transfers)

Data for analysis tasks — MDS disk, 30 TB

100 MB/sec (file transfers)

Interesting Data — HPSS Robot, 300 TB

All Data — HPSS Shelf, 3 PB

Bulk bandwidth numbers meet estimated requirements assuming 100% efficiency.

How to achieve bulk bandwidth?

What fraction of data transfered is useful to programs?!!!

Keep recently accessed data on disk, but manage it so unused data does not waste space.

Try to arrange that 90% of file access is to disk and only 10% are retrieved from tape.

# Solution Components

- **Split event into components across different files so that most bytes read are used**
  - Raw, tracks, hits, tags, summary, trigger, …
- **Optimize file size so tape bandwidth is not wasted**
  - 1GB files, → means different # of events in each file
- **Coordinate file usage so tape access is shared**
  - Users select all files at once
  - System optimizes retrieval and order of processing
- **Use disk space & bandwidth efficiently**
  - Operate disk as cache in front of tape

# STAR Event Model



T. Ullrich, Jan. 2000

# Analysis of Events

- **1M events = 100GB – 1TB**

  – **100 – 1000 files (or more if not optimized)**

- **Need to coordinate event associations across files**

- **Probably have filtered some % of events**

  – **Suppose 25% failed cuts after trigger selection**

    - **Increase speed by not reading these 25%**

- **Run several batch jobs for same analysis in parallel to increase throughput**

- **Start processing with files already on disk without waiting for staging from HPSS**
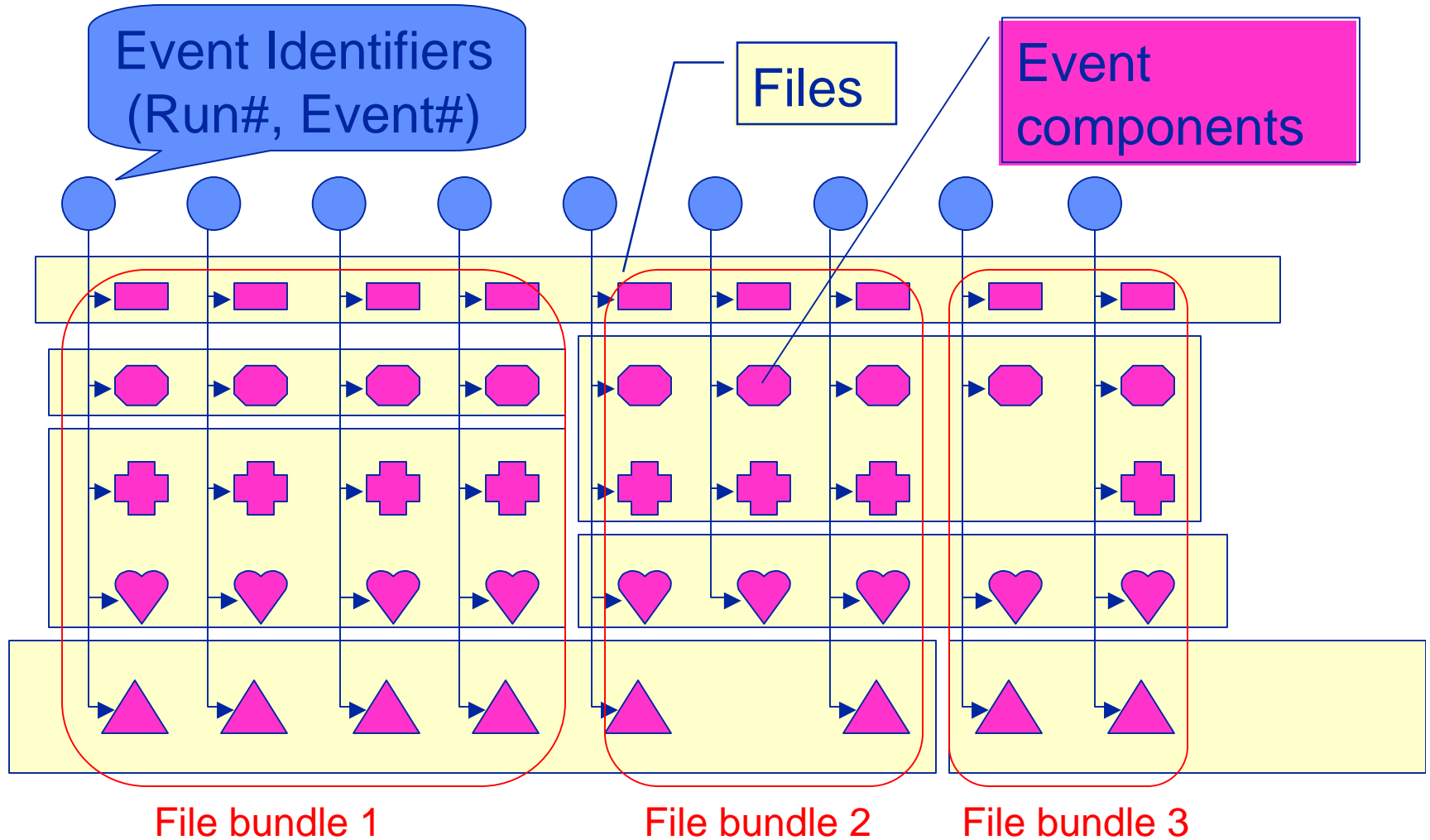
# GCA System Overview

*fileCatalog, tagDB* and  GCA

# STACS: STorage Access Coordination System

# Organization of Events in Files



Event Identifiers (Run#, Event#)

Files

Event components
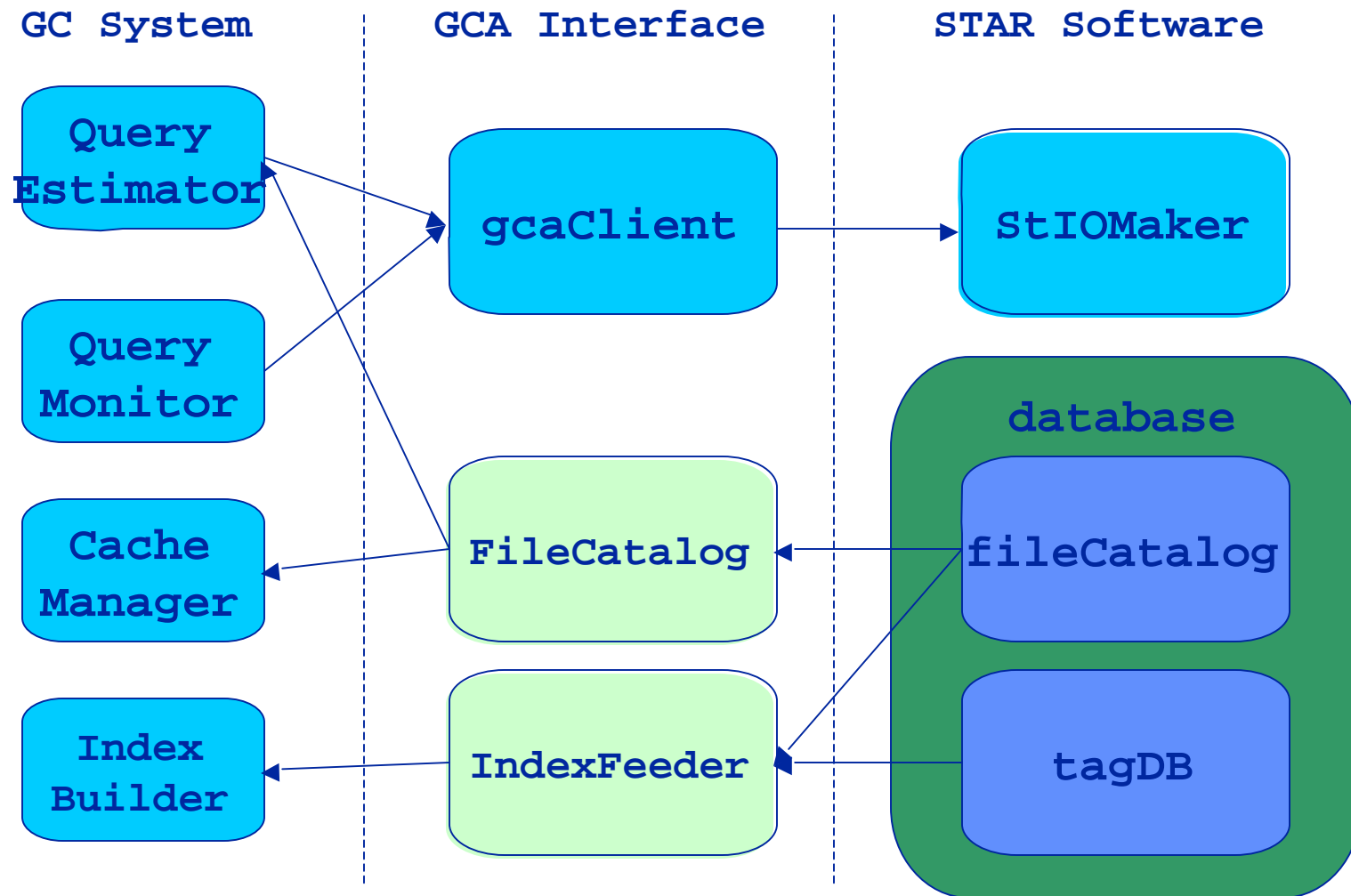
File bundle 1

File bundle 2

File bundle 3

# The Details

- **Range-query language, or query by event list**
  - **"NLa>700 && run=101007",**
  - **{e1,r27012;e3,r27014;e7;r27017 …}**
  - **Select components: dst, geant, …**
- **Query estimation**
  - **# events, # files, # files on disk, how long, …**
  - **Avoid executing incorrect queries**
- **Order optimization**
  - **Order of events you get maximizes file sharing and minimizes reads from HPSS**
- **Policies**
  - **# of pre-fetch, # queries/user, # active pftp connections, …**
  - **Tune behavior & performance**
- **Parallel processing**
  - **Submitting same query token in several jobs will cause each job to process part of that query**

# Interfacing GCA to STAR

**GC System**　　　　**GCA Interface**　　　　**STAR Software**

| Query Estimator |
| Query Monitor |
| Cache Manager |
| Index Builder |

| gcaClient |
| FileCatalog |
| IndexFeeder |

| StIOMaker |

**database**
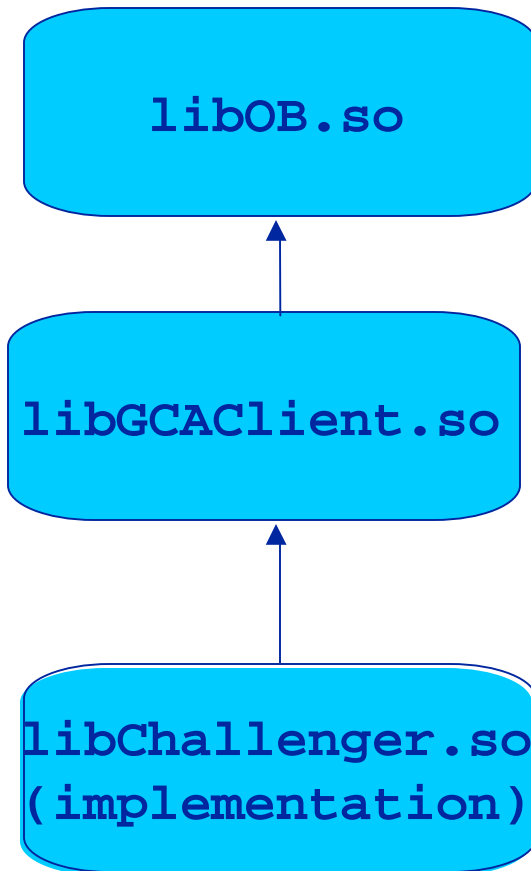
| fileCatalog |
| tagDB |

# Limiting Dependencies

### STAR-specific  &  GCA-dependent

- # IndexFeeder server
  - – **IndexFeeder read the "tag database" so that GCA "index builder" can create index**

- # FileCatalog server
  - – **FileCatalog queries the "file catalog" database of the experiment to translate fileID to HPSS & disk path**

- # gcaClient interface
  - – **Experiment sends queries and get back filenames through the gcaClient library calls**

# Eliminating Dependencies

**CORBA + GCA software**

**ROOT + STAR Software**

```
libOB.so
```

```
libGCAClient.so
```

```
libChallenger.so
(implementation)
```

```
TNamed
```

```
<<Interface>>
StFileI
```

```
StChallenger
::Challenge()
```

```
StIOMaker
```

# STAR *fileCatalog*

- **Database of information for files in experiment. File information is added to DB as files are created.**

- **Source of File information**

  – **for the experiment**

  – **for the GCA components (Index, gcaClient,...)**

# Cataloguing Analysis Workflow



Job configuration manager

fileCatalog

Job monitoring system

*fileCatalog, tagDB* and GCA

# Transactionless Solution

- **MySQL:**
  - **no views**
    - reader has to join tables
  - **no transactions**
    - db snapshot during the update may be inconsistent
    - updates may be long (= no table locking for read)
- Experiment: few writers, hundreds of readers
- Compromise:
  - **use pre-calculated views (= extra tables)**
    - **fileCatalog data are duplicated**
  - **views update is quick (table can be locked)**
    - **clients read see consistent data at any time**

# Problem:  SELECT NLa>700

**ntuple**

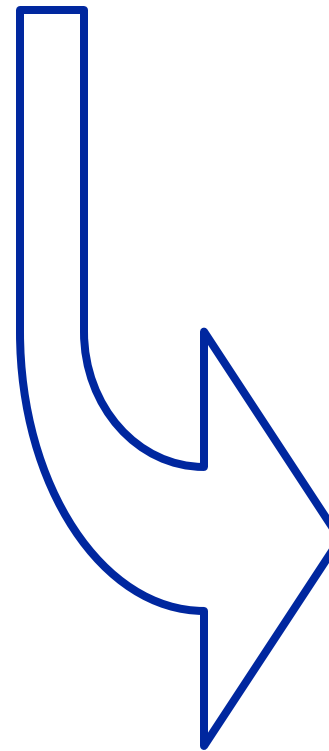**read all events**

| Event # | NLa |
|---------|-----|
| 1 | 731 |
| 2 | 800 |
| 3 | 345 |
| 4 | 543 |
| 5 | 567 |

**index**

**read selected events**

| NLa | Event # |
|-----|---------|
| 345 | 3 |
| 543 | 4 |
| 567 | 5 |
| 731 | 1 |
| 800 | 2 |

# STAR Tag Structure Definition

./pams/global/idl/FlowTag.idl

Version: [ .DEV ] [ DEV00 ] [ SL00b_2 ] [ SL99f ]

```
1  //
2  // $Id: FlowTag.idl,v 1.3 2000/01/13 23:18:06 snelling Exp $
3  //
4  // Event by event flow tag
5  //
6  // $Log: FlowTag.idl,v $
7  // Revision 1.3  2000/01/13 23:18:06   snelling
8  // Changed sum pt to mean pt
9  //
10 // Revision 1.2  1999/11/16 20:59:40   snelling
11 // Removed unused tags and added 6th harmonic
12 //
13 // Revision 1.1  1999/02/09 21:42:21   wenaus
14 // Final (?) versions of MDC2 PWG tags
15 //
16 // The tags are defined for 4 subevents (a,b,c,d) and 6 harmonics
17
18 struct FlowTag {
19   float qxa[6], qxb[6], qxc[6], qxd[6];       /* x component Q vector */
20   float qya[6], qyb[6], qyc[6], qyd[6];       /* y component Q vector */
21   long  na[6], nb[6], nc[6], nd[6];           /* multiplicity */
22   float mpta[6], mptb[6], mptc[6], mptd[6];   /* mean pt */
23 };
```

**Selections like**

$\sqrt{qxa^2 + qxb^2} > 0.5$

**can not use index**

# STAR Tag Database Access

# TagDB in ROOT Files

- **Tag data are stored in ROOT TTree files**

- **Branches (3 out of 6 requested) saved in the split mode**
  - StrangeTag
  - FlowTag
  - ScaTag

- **173 physics tags [int/float] out of 500 requested**

- **Disk resident tag files name+path are stored in the MySQL *fileCatalog* database**

- **Files are selected from database in the ROOT CINT macro through the ROOT-MySQL interface and are chained for further selections by user**

# MDC3 Index

- **6 event components:**
  - fzd
  - geant
  - dst
  - tags
  - runco
  - hist

- **179 physics tags:**
  - StrangeTag
  - FlowTag
  - ScaTag

- **120K events**

- **8K files**

# GCA MDC3 Integration Workshop

**http://www-rnc.lbl.gov/GC/meetings/14mar00/default.htm**

**14-15 March 2000**

**Goals:**

| status | goal | description / summary (as of 16Mar2000) |
|--------|------|------------------------------------------|
| done | 1 | **Build index on new STAR files** |
| | | The index was build (several times) on the new STAR MDC3 data. This consisted of about 5,000 events. By the end of next week (start of MDC3) STAR expects about 140K events to put in the GC index. Sasha is continuing to accumulate additional event tag files as they are available. |
| done | 2 | **Check that GCAClient and MinimalQuery work** |
| | | Modifications to GCAClient and the MinimalQuery (& MinimalQuery1) test programs were completed for the updated version of STACS, including the new file bundle flag on the iterator. |
| done | 3 | **Run MinimalQuery on linux** |
| | | GCAClient & test program was compiled, run successfully on linux as well as Solaris. This included modifications to the Makefile to build both on linux & solaris. |
| done | 4 | **Run multiple MinimalQuery simultaneously** |
| | | Run on linux. Not verified yet on solaris. |
| done | 5 | **Test index update** |
| | | The feature of being able to update (add new events) to an existing index was justed added. This feature was first tested during this period. A number of bug fixes were made and the basic procedure is working. John is continuing to investigate one or two bugs before the procedure is declared reliable. |
| done | 6 | **Test index update while queries are running** |
| | | This is a functionality test and was successful. Any remaining work on the udpate functionality is not related to interlocks with running queries. |
| done | 6.1 | **update between queries** |
| | | This check is to run a query before the update and then after and verify that the results are accurate. This was successful. |
| done | 6.2 | **update while new queries are being submitted** |
| | | This tests the interlock mechanism so that queries do not run during the update process. This was successful. |
| in progress | 7 | **Integrate GCAClient into root4star** |
| | | This is the final work to connect the GCA to STAR data analysis. There were various discussions among Victor, Sasha, Jeff, Frank, Dave, Doug. The basic idea of how to incorporate the GCAClient into StIOMaker has been worked out by Victor, Sasha & Jeff. Sasha & Victor will work on it. |

# User Query

## ROOT Session:



```
rcas6023:/star/u2c/vanyashi/gc/StGCAClient
1 mBeamPolarizationWest_0
1 mBeamPolarizationWest_1
1 mBeamPolarizationWest_2
1 mBImpact
1 mPhImpact
0 mGenerType
0 mBunchCrossingNumber
0 mEventNumber
0 mEventTime
0 mEventDate
0 mProdTime
0 mProdDate
qM: 0x86cb588
qE: 0x86cb158
fC: 0x86cba40
qF: 0x86cbd48
a->Init()
 *** OidSource is not set.

   Submitting query: SELECT dst
                     WHERE  -5<=qxa_3<0.3 && 22>qxc

qoF:: query created
qoF:: query added to list
query 0x86d0ce4
Full estimate is 205 events in 161 files ( unknown MBs).
```

```
rcas6023:/star/u2c/vanyashi/gc/StGCAClient
root.exe [0]
Processing test.C...
StGCAAdapter::LoadGCAServer: libStGCAClient.so loaded
StGCAAdapter::LoadGCAServer: new StGCAServer created
          StGCAServer::Init messages:
          I will not attempt to follow refs returned via the iterator.
gcaResources:  Attempting to read configFile /star/rcf/GC/MDC3/stacs.rc
Using configuration file "/star/rcf/GC/MDC3/stacs.rc".
Narrowing QE reference found in /star/rcf/GC/MDC3/logs/SM_QE.ref
Converting (string_to_object) IOR:00000000000001549444c3a736d457374696c
00000005c00010000000000137273756e30302e7263662e626e6c2e676f76000006be000
2e676f763a5175657279457374696d6174f723a303a3a49523a736d457374696d61746f
Converted string_to_object
returning from findObjViaStringFile...
A Query Estimator has been contacted.
Converting (string_to_object) IOR:00000000000001849444c3a716d4576656e74
00000005c00010000000000137273756e30302e7263662e626e6c2e676f76000006c1000
2e676f763a51756572794d6f6e69746f723a313a3a49523a716d4576656e744974657261
Converted string_to_object
returning from findObjViaStringFile...
A Query Monitor is available to your OrderOptIterator.
Narrowing FileCatalog reference found in /star/rcf/gc/GCdev/FC/FileCatal
Converting (string_to_object) IOR:00000000000001449444c3a46696c65436174
40001000000000137273756e30302e7263662e626e6c2e676f760000883b00000000001
Converted string_to_object
returning from findObjViaStringFile...
A File Catalog has been found.
You are connected to a Query Factory.
Index Information
Name = simulated data for MDC3
Description =51749 events, 179 attributes, 6 components (all NULL FIDs a
00
number of components = 6
dst
fzd
geant
hist
runco
tags
0 tags
```

# Conclusion

- **GCA developed a system for optimized access to multi-component event data files stored in HPSS**

- **General CORBA interfaces are defined for interfacing with the experiment**

- **A client component encapsulates interaction with the servers and provides an ODMG-style iterator**

- **Has been tested up to 10M events, 7 event components, 250 concurrent queries**

- **Has been integrated with the STAR experiment ROOT-based I/O analysis system**